

Ora devo verificare che effettivamente p_k è una direzione di decrescita:

$$(\nabla \phi(x))^T p_k = -r_k^T (r_k + \beta_k p_{k-1}) =$$

$$\stackrel{!}{=} -r_k^T r_k - \beta_k r_k^T p_{k-1} =$$

Ora dimostreremo che β_k è uguale a zero e quindi (viene per come abbiamo scelto α_k , non dipende nemmeno da p_k)

$$\stackrel{!}{=} -\|r_k\|_2^2 < 0 \quad \alpha_k$$

PROP: 1) $\forall k \quad r_{k+1}^T p_k = 0$ (vale in generale)

2) $\forall k \quad r_{k+1}^T r_k = 0$ quindi le coppie di residui sono ortogonali (in senso euclideo) (vale solo per CG)

$$\text{OSS: } r_{k+1} := b - Ax_{k+1} = b - A(x_k + \alpha_k p_k) = (b - Ax_k) - \alpha_k A p_k = \\ \stackrel{!}{=} r_k - \alpha_k A p_k$$

$$\text{dim: (1) } r_{k+1}^T p_k = (r_k - \alpha_k A p_k)^T p_k = r_k^T p_k - \alpha_k p_k^T A p_k \stackrel{!}{=} 0$$

$$\text{def di } \alpha_k = \frac{r_k^T p_k}{p_k^T A p_k}$$

$$(2) \quad r_{k+1}^T r_k = r_{k+1}^T (p_k - \beta_k p_{k-1}) = r_{k+1}^T p_k - \beta_k r_{k+1}^T p_{k-1} = \\ \stackrel{!}{=} -\beta_k (r_k - \alpha_k A p_k)^T p_{k-1} = -\beta_k r_k^T p_{k-1} + \alpha_k \beta_k p_k^T A p_{k-1} = 0$$

perché $p_k^T A p_{k-1}$ per def dei p_k .

Ora vogliamo vedere l'algoritmo del CG. Prima di scrivere l'algoritmo vediamo delle RISCRIURE di α_k e β_k per RIDURRE il COSTO COMPUTAZIONALE. Abbiamo

$$\alpha_k := \frac{r_k^T p_k}{p_k^T A p_k} = \frac{r_k^T r_k}{p_k^T A p_k}$$

NB non è che $p_k = r_k$ ma quando facciamo il prodotto con r_k^T vengono uguali

$$\beta_k := -\frac{r_k^T A p_{k-1}}{p_{k-1}^T A p_{k-1}} = -\frac{r_k^T r_{k-1}}{r_{k-1}^T r_{k-1}}$$

← sono spariti i prodotti matriciali che hanno costo computazionale maggiore. In più $r_k^T r_{k-1}$ l'ho già calcolato al passo precedente.

$$\text{dim: } r_k^T p_k = r_k^T (r_k + \beta_k p_{k-1}) = r_k^T r_k + \beta_k r_k^T p_{k-1} \stackrel{!}{=} r_k^T r_k$$

Vale in generale, perché (1) vale in generale, mi piace perché sta riscrivendo α_k , che non "dipende" dal metodo CG

$$\bullet p_k^T r_{k-1} = p_k^T (r_k + \alpha_{k-1} A p_{k-1}) = p_k^T r_k + \alpha_{k-1} p_k^T A p_{k-1} = \\ \stackrel{!}{=} p_k^T r_k = r_k^T r_{k-1}$$

Ma ho anche:

$$p_k^T r_{k-1} = (r_k + \beta_k p_{k-1})^T r_{k-1} = r_k^T r_{k-1} + \beta_k p_{k-1}^T r_{k-1} = \beta_k r_{k-1}^T r_{k-1}$$

Uguagliando le due quantità si ricava β_k .

ALGORITMO CG: pcg di Matlab, lavora in due modi.

Dato x_0 (in genere $x_0 = 0$), facciamo

$$1. \quad k=0, \quad r_0 = b - Ax_0$$

→ nuova ipotesi non dovuto non $r_0 = 0$

2. Test d'arresto se " $r_k = 0$ " stop

$$3. \quad \beta_k = \frac{r_k^T r_k}{r_{k-1}^T r_{k-1}} \quad (\beta_0 = 0 \text{ se } k=0)$$

$$p_k = r_k + \beta_k p_{k-1}$$

$$\alpha_k = \frac{r_k^T r_k}{p_k^T A p_k}$$

→ A, b, m, x
(event M precond)
→ $A, b, (M) +$
function per fare i 2 pt.

• Calcola Ax
con prodotto
matr. veloce
• $Mx = r$
con LU

21/11/2012

(12)

$$x_{k+1} = x_k + \alpha_k p_k$$

$$r_{k+1} = r_k - \alpha_k A p_k \quad (= b - A x_{k+1} \text{ ma nella pratica si usa l'altra})$$

4. $k = k+1$ e vai a 2.

Il costo pesante del metodo è il prodotto matrice vettore $A p_k$ che mi comporre 2 volte (ma è sempre lo stesso) per come ho ottimizzato scrivendo α_k e β_k come ho fatto. $\Rightarrow O(n^2)$

Osservo che se avessi calcolato $r_{k+1} = b - A x_{k+1}$ avrei avuto 2 prodotti matrice vettore diversi!

Poi ho dei prodotti scalari tra vettori che però costano meno (e che tra l'altro si ripetono uguali).

Quindi il costo computazionale è un $O(n^2)$

OSS: È vero che in teoria se calcolo $r_{k+1} = b - A x_{k+1}$ o con $r_k - \alpha_k A p_k$ è uguale ma in pratica con l'aritmetica floating point può non essere uguale.

In pratica quando $k \approx \sqrt{n}$ se non ho ancora raggiunto la condizione d'arresto calcolo per una volta $r_{k+1} = b - A x_{k+1}$ e poi riprendo a calcolare nel solito modo. Così non incremento il costo comp. perché per una volta che lo faccio cambio poco, ma recupero precisione.

Il test d'arresto che si fa è quello del residuo relativo:

$$\frac{\|r_k\|_2}{\|b\|_2} = \frac{\|b - A x_k\|_2}{\|(b - A x_0)_{x=0}\|_2} < \varepsilon$$

($\|b\|_2$ lo calcolo una volta sola e quindi non incide sul costo)

Ora vediamo una proposizione che ci fa passare dall'ortogonalità a coppie dei residui, all'ortogonalità tra tutti, in modo che così ho dimostrato che CG si conclude al più in n passi.

PROP: Sia $r_0 \neq 0$, $s \geq 1$ t.c. $r_k \neq 0 \quad \forall k \leq s$ (i.e. non ho già trovato la soluzione esatta)

$$\Rightarrow \forall k, j = 0, \dots, s, \quad k \neq j \quad \begin{aligned} r_k &\perp r_j \\ p_k &\perp_A p_j \end{aligned}$$

dim.: Per induzione. • Prendo $s=1$ $k=1$ e $j=0$

$$r_1^T r_0 = (r_0 - \alpha_0 A p_0)^T r_0 = r_0^T r_0 - \alpha_0 p_0^T A r_0 = 0 \quad (\text{ma qui lo so già per la prop precedente})$$

$$p_1^T A p_0 = 0 \quad \text{per costruzione.}$$

• Sia ora $s > 1$. hp di induzione: $\forall k, j = 0, \dots, s, \quad k \neq j \quad r_k^T r_j = 0$ e $p_k^T A p_j = 0$

Voglio dimostrare che:

$$\forall j = 0, \dots, s-1 \quad r_{s+1}^T r_j = 0 \quad \text{e} \quad p_{s+1}^T A p_j = 0$$

↳ tutto per $j=s$ è a posto per la PROP precedente.

$$\begin{aligned} r_{s+1}^T r_j &= (r_s - \alpha_s A p_s)^T r_j = r_s^T r_j - \alpha_s p_s^T A r_j = \\ &= -\alpha_s p_s^T A (p_j - \beta_j p_{j-1}) = -\alpha_s p_s^T A p_j + \alpha_s \beta_j p_s^T A p_{j-1} = 0 \quad \text{per hp di induzione} \end{aligned}$$

$$p_{s+1}^T A p_j = (r_{s+1} + \beta_{s+1} p_s)^T A p_j = r_{s+1}^T A p_j + \beta_{s+1} p_s^T A p_j = 0 \quad \text{per hp di ind.}$$

Dato che ho $r_{k+1} = r_k - \alpha_k A p_k$ posso scrivere $A p_j = \alpha_j^{-1} (r_j - r_{j+1})$

$$= \alpha_j^{-1} (r_{s+1}^T r_j - r_{s+1}^T r_{j+1}) = 0 \quad \text{per la tesi già dimostrata.}$$

TEO: Il CG applicato a $Ax=b$, $A \in \mathbb{R}^{n \times n}$ semidef. positiva, genera la soluzione esatta in al più n passi. (in aritmetica esatta)

dim: \exists solo n vettori \perp tra loro. Ci vuole di meno se la soluzione vive in un sottospazio di dim $<$

OSS: Abbiamo visto che se $x^* = A^{-1}b$ allora $\phi(x^*+e) = \phi(x^*) + \frac{1}{2} e^T A e$

se scriviamo $x = x^* + e \rightarrow e = x - x^*$ possiamo riscrivere la relazione come:

$$\phi(x) = \phi(x^*) + \frac{1}{2} (x - x^*)^T A (x - x^*)$$

Minimizzare $\phi(x_k)$ equivale a minimizzare $\|e_k\|_A = \|A^{1/2} e_k\|_2$ o.e. $(A^{1/2} = Q \Lambda^{1/2} Q^H)$ dove $Q \in \mathbb{C}^n$ sono le matrici della fatt. $A = Q \Lambda Q^H$ agli autovalori.)

Quindi il metodo funziona anche se $A \in \mathbb{C}^{n \times n}$ se lo vediamo in questa chiave di lettura.

TEO: L'errore al passo $k+1$ è A -ortogonale a tutte le direzioni di ricerca precedenti:

$$i.e. e_{k+1}^H A p_j = 0 \quad \forall j = 0, \dots, k$$

Inoltre: $\text{span}\langle p_0, p_1, \dots, p_k \rangle = \text{span}\langle A e_0, \dots, A^{k+1} e_0 \rangle$ (è SPAZIO di KRYLOV)

E quindi e_{k+1} è il vettore di norma A più piccola nello spazio $e_0 + \text{span}\langle A e_0, \dots, A^{k+1} e_0 \rangle$. (forse lo dimostra la prossima volta o forse no)

Vediamo ora le STIME di VELOCITA' di CONVERGENZA del CG

TEO: $\frac{\|e_k\|_A}{\|e_0\|_A} \leq \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \left\{ \max_{i=1, \dots, n} |p_k(\lambda_i(A))| \right\}$

(autovalori di A)
polinomi monici di grado k

N.B.: Dipende da tutti gli autovalori della matrice A !

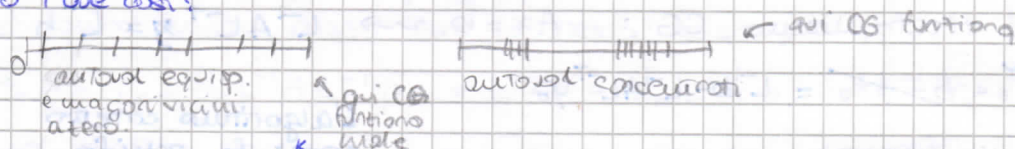
Dovremmo aver già visto il teo:

TEO: $\frac{\|e_k\|_A}{\|e_0\|_A} \leq 2 \left(\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k$ con $\kappa_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$

la dimostrazione di questo risultato fa pesante uso di teoria dell'approssimazione (non è propriamente algebrica), e ci è fuorviante, perché ci fa vedere solo la dipendenza da λ_{\max} e λ_{\min} .

Questa seconda stima è pessimistica (e in pratica non mi dice nulla, perché dice se A è mal condizionata potrebbe dare problemi.)

Il nostro Teorema invece è più fine, per esempio fa differenza tra i due casi:



dim: $e_k = e_0 + \sum_{i=1}^k \chi_i A^i e_0 = (1 + \sum_{i=1}^k \chi_i A^i) e_0 = p_k(A) e_0$

(TEO) con $p_k(x) = 1 + \sum_{i=1}^k \chi_i x^i \in \mathcal{P}_k$ oss: $p_k(0) = 1$!

Scrivo $A = Q \Lambda Q^H$

$\|e_k\|_A = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|A^{1/2} p_k(A) e_0\|_2 = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|p_k(A) A^{1/2} e_0\|_2 = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|Q p_k(\Lambda) Q^H A^{1/2} e_0\|_2$

(posso far transitare la matrice a dx)

Osserviamo che il minimo dipende da e_0 e questo ci dà un po' noia, quindi facciamo una maggiorazione.

Si può mostrare che questa $\epsilon \leq \min \|Q p_c(\lambda) Q^H\|_2 \|A^{1/2} e_0\|_2 \stackrel{Q \text{ unitaria}}{=} \min \|p_c(\lambda)\| \|e_0\|_A =$
 maggiore non è troppo = $\min \{ \max_{i=1, \dots, n} |x_i(A)| \} \|e_0\|_A$
 poco fine (ci sono morici per cui vale l'uguale).

Vediamo come ci comportiamo nel caso in cui la situazione sia brutta:

PRECONDIZIONAMENTO:

- (1) $Ax = b \iff M^{-1}Ax = M^{-1}b$ con M opportuna (precondizionamento a Sx)
 Σ_A spettro di A (brutto) Voglio che lo spettro di $M^{-1}A$ sia migliore per il metodo (di Krylov in generale).

OSS: Anche se in teoria con il precondizionamento mi porta a un problema ben condizionato, in pratica non sarà mai così, perché se A è mal condizionato non calcolerò bene il prodotto matrice vettore ($M^{-1}A$).

Quindi da un problema mal condizionato spero di ottenere un problema sempre mal condizionato ma che ha velocità di convergenza maggiore!

23/11/2012

- (2) Analogamente è possibile fare il precondizionamento a dx :

$$Ax = b \iff AM^{-1}\hat{x} = b \quad \text{con } x = M^{-1}\hat{x}$$

- (3) Per il CG si usa solitamente un precondizionamento simmetrico:

Scriviamo la fattorizzazione di Cholesky: $M = LL^T$ (uso un precond. sdP)

$$Ax = b \iff L^{-1}AL^{-T}y = L^{-1}b \quad \text{con } x = L^{-T}y$$

Come scegliamo M ? Tipicamente si usa

$M =$ fattorizzazione incompleta di Cholesky \tilde{L}

Perché vorrei una matrice "bella" della forma di A , detto che nei nostri problemi A è tridiagonale a blocchi vogliamo M e L fatte così. La fatt. di Cholesky riempie la matrice, l'incompleta invece lo fa come se fosse completa, ma lascia zero dove è zero.

M si sceglie in modo che $M^{-1}A$ (o AM^{-1} o $L^{-1}AL^{-T}$) "approssimi" l'identità. In che senso approssima dipende dal metodo.

Inoltre chiaramente M deve essere t.c. risolvere $Mz = r$ sia più facile che risolvere $Ax = b$.

Vediamo come trasformiamo CG: $Ax = b \iff L^{-1}AL^{-T}y = L^{-1}b$

1. $k=0 \quad \tilde{r}_0 = b - Ax_0 = L^{-1}b - L^{-1}AL^{-T}y_0$

2. Se $\tilde{r}_k = 0$ mi fermo

3. $\tilde{\beta}_k = \frac{\tilde{r}_k^T \tilde{r}_k}{\tilde{r}_k^T \tilde{r}_k} \quad \tilde{\beta}_0 = 0$

$$\tilde{p}_k = \tilde{r}_k + \tilde{\beta}_k \tilde{p}_{k-1}$$

$$\tilde{\alpha}_k = \frac{\tilde{r}_k^T \tilde{r}_k}{\tilde{p}_k^T A \tilde{p}_k} \quad \tilde{p}_k^T L^{-1}AL^{-T} \tilde{p}_k$$

$$x_{k+1} = x_k + \tilde{\alpha}_k \tilde{p}_k \rightarrow y_{k+1} = y_k + \tilde{\alpha}_k \tilde{p}_k$$

$$\tilde{r}_{k+1} = \tilde{r}_k - \tilde{\alpha}_k A \tilde{p}_k$$

L'algoritmo corretto così fa schifo, ora vediamo effettivamente come si fa.

ALGORITMO PCG: su $Ax = b$

1. $k=0$ $r_0 = b - Ax_0$
2. se $r_k = 0$ mi fermo
3. calcolo z_k t.c. $Mz_k = r_k$

$$\tilde{\beta}_k = \frac{z_k^T r_k}{z_k^T r_{k-1}} \quad \text{* già calcolato al passo precedente}$$

$$p_k = z_k + \tilde{\beta}_k p_{k-1}$$

$$\tilde{\alpha}_k = \frac{z_k^T r_k}{p_k^T A p_k}$$

$$x_{k+1} = x_k + \tilde{\alpha}_k p_k$$

$$r_{k+1} = r_k - \tilde{\alpha}_k A p_k$$

$k = k+1$ e torno a 2.

← L'unico aumento di costo computazionale sta qui, ovvero nella risoluzione di un sistema lineare in M ad ogni passo. (il resto costa uguale)

N.B.: L'algoritmo è davvero il CG applicato al sistema preconditionato, è solo scritto in modo da diminuire il costo computazionale (ottimizzare)

Vediamolo:

$$\tilde{r}_k = L^{-1}b - L^{-1}AL^{-T}y_k = L^{-1}(b - AL^{-T}y_k) \stackrel{x_k = L^{-T}y_k}{=} L^{-1}(b - Ax_k) = L^{-1}r_k$$

$$\tilde{\beta}_k = \frac{\tilde{r}_k^T \tilde{r}_k}{\tilde{r}_{k-1}^T \tilde{r}_{k-1}} = \frac{r_k^T L^{-T} L^{-1} r_k}{r_{k-1}^T L^{-T} L^{-1} r_{k-1}} \stackrel{M=LL^T}{=} \frac{r_k^T M^{-1} r_k}{r_{k-1}^T M^{-1} r_{k-1}} \stackrel{z_k = M^{-1} r_k}{=} \frac{z_k^T r_k}{z_{k-1}^T r_{k-1}}$$

Rimane da vedere che possiamo usare p_k al posto di \tilde{p}_k , ecc..

Abbiamo che $y = L^T x$ e quindi $\tilde{p}_k = L^T p_k$.

$$p_k = L^{-T} \tilde{p}_k = L^{-T}(\tilde{r}_k + \tilde{\beta}_k \tilde{p}_{k-1}) = L^{-T}(L^{-1}r_k + \tilde{\beta}_k L^T p_{k-1}) = L^{-T} L^{-1} r_k + \tilde{\beta}_k p_{k-1} = M^{-1} r_k + \tilde{\beta}_k p_{k-1} = z_k + \tilde{\beta}_k p_{k-1} \quad \text{ok.}$$

$$\tilde{\alpha}_k = \frac{\tilde{r}_k^T \tilde{r}_k}{\tilde{p}_k^T L^{-1} A L^{-T} \tilde{p}_k} \stackrel{\text{visto per } \beta}{=} \frac{z_k^T r_k}{p_k^T A p_k} \rightarrow \text{per (*)}$$

$$y_{k+1} = y_k + \tilde{\alpha}_k \tilde{p}_k = L^T x_k + \tilde{\alpha}_k L^T p_k = L^T(x_k + \tilde{\alpha}_k p_k)$$

Ma ovviamente $y_{k+1} = L^T x_{k+1}$ e quindi uguagliando le due espressioni dtego:

$$x_{k+1} = x_k + \tilde{\alpha}_k p_k$$

e tanto quanto posso fare la stessa cosa per r_{k+1} :

$$\tilde{r}_{k+1} = \tilde{r}_k - \tilde{\alpha}_k L^{-1} A L^{-T} \tilde{p}_k = L^{-1} r_k - \tilde{\alpha}_k L^{-1} A p_k = L^{-1}(r_k - \tilde{\alpha}_k A p_k)$$

ma d'altra parte $\tilde{r}_{k+1} = L^{-1} r_{k+1}$, da cui $r_{k+1} = r_k - \tilde{\alpha}_k A p_k$.

OSS: il costo del PCG è tipo: prodotto matrice-vettore + risoluzione sistema lineare in M ($Mz_k = r_k$)
con A ($A p_k$)
= costo di CG

Il punto è, se uso CG ho $O(n^2)$ come costo comp., se davvero mi tocca fare n iterazioni arrivo a $O(n^3)$ ed è inutile, perché potrei usare Gauss!

Se invece la matrice è sparsa CG mi costa $O(n)$ e anche la risoluzione di $Mz=r$ è $O(n)$ (se uso Cholesky incompleta), se ho ancora n iter. arrivo a $O(n^2)$ che è ancora uguale a Gauss per matrici sparse!

Quindi mi fa gioco sul numero di iterazioni! il preconditionatore può fare davvero molto.

Vediamo ulteriori risultati sulla convergenza del metodo del gradiente coniugato. (basati su AXELSSON-LINDSCOG '86)

Consideriamo l'intervallo $[a, b]$, $a = \lambda_{\min}$, $b = \lambda_{\max}$ allora

$$p_k^*(\lambda) = T_k\left(\frac{b+a-2\lambda}{b-a}\right) / T_k\left(\frac{b+a}{b-a}\right)$$

→ trasform. affine $[-1, 1] \rightarrow [a, b]$

ovvero $p_k^*(\lambda)$ è il polinomio di miglior approssimazione del TEO di 2 lezioni fa (min max...) e T_k è il polinomio di Chebyshev di grado k :

$$T_k(x) = \cos(k\theta) \quad \text{con } \theta = \arccos x$$

$$= \frac{1}{2} \left[(x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right] \quad \text{su } [-1, 1].$$

Abbiamo:

$$\max_{\lambda \in [a, b]} |p_k^*(\lambda)| = \frac{1}{T_k\left(\frac{b+a}{b-a}\right)} = \frac{2\sigma^k}{1+\sigma^{2k}} \quad \text{con } \sigma = \frac{1 - \sqrt{\frac{a}{b}}}{1 + \sqrt{\frac{a}{b}}} = \frac{\sqrt{b} - \sqrt{a}}{\sqrt{b} + \sqrt{a}}$$

Quindi se voglio che l'errore relativo sia $< \varepsilon$ (i.e. $\frac{\|e_k\|_A}{\|e_0\|_A} < \varepsilon$) chiedo che

$$\frac{2\sigma^k}{1+\sigma^{2k}} \leq \varepsilon$$

che facendo i conti diventa:

$$k \geq \left\lceil \frac{\log\left(\frac{1}{\varepsilon} + \sqrt{\frac{1}{\varepsilon^2} - 1}\right)}{\log(\sigma^{-1})} \right\rceil$$

→ si può trascurare perché se ε è piccolo $\frac{1}{\varepsilon^2}$ è grande

$$\leadsto k \geq \left\lceil \frac{\log\left(\frac{2}{\varepsilon}\right)}{\log(\sigma^{-1})} \right\rceil = k^*(a, b, \varepsilon)$$

OSS: la prima cosa che noto è che la crescita rispetto a ε è logaritmica, e quindi ci piace.

Ora dobbiamo vedere cosa succede rispetto a σ^{-1} . Abbiamo:

$$\sigma^{-1} = \frac{\sqrt{b} + \sqrt{a}}{\sqrt{b} - \sqrt{a}} = \frac{(1 + \sqrt{\frac{b}{a}})}{-(1 - \sqrt{\frac{b}{a}})} = \frac{\sqrt{K_2(A)} + 1}{\sqrt{K_2(A)} - 1} = 1 + \frac{2}{\sqrt{K_2(A)} - 1}$$

Pensiamo di avere una successione di sistemi lineari $A_n x_n = b_n$ (che poi è quello che abbiamo quando per esempio vediamo gli elem. finiti e facciamo tendere $h \rightarrow 0$), vogliamo vedere cosa succede per $n \rightarrow \infty$

Consideriamo per esempio gli elem. finiti, abbiamo che

$$K_2(A_n) \nearrow \infty \quad \text{per } n \rightarrow \infty$$

$$\text{e quindi } \sigma^{-1} = 1 + \varepsilon_n \quad \text{con } \varepsilon_n = \frac{2}{\sqrt{K_2(A_n)} - 1} \searrow 0$$

Dunque: $\log(\sigma^{-1}) = \log(1 + \varepsilon_n) = \varepsilon_n + O(\varepsilon_n^2)$ e in conclusione

$$\begin{aligned} k^*(a, b, \varepsilon) &= \left\lceil \frac{\log\left(\frac{2}{\varepsilon}\right)}{\varepsilon_n + O(\varepsilon_n^2)} \right\rceil = \\ &= \left\lceil \log\left(\frac{2}{\varepsilon}\right) \varepsilon_n^{-1} (1 + O(\varepsilon_n)) \right\rceil = \\ &= \left\lceil \log\left(\frac{2}{\varepsilon}\right) \frac{\sqrt{K_2(A_n)} + 1}{2} (1 + O(\varepsilon_n)) \right\rceil \end{aligned}$$

Quindi moralmente stiamo dicendo che se $K_2(A)$ è alto ci vogliono più iterazioni, o ancora meglio se $[a, b]$ è grande ci vogliono più iterazioni. (ci piacciono gli autov. tutti concentrati e lontani da zero)

CONSIDERAZIONI: Tipicamente la matrice che abbiamo noi (FEM/DF)
FILOSOFIA: avrà autovalori più o meno equispaziati:



Il punto è trovare un preconditionatore che ci concentri gli autovalori e li porti in una configurazione tipo:



E gli autovalori isolati? Quelli in pratica li trascureremo usando un condizionamento essenziale (che tiene conto degli autovalori concentrati) e poi mi occupo solo dopo degli outlier.

VARIANZI SUL TEMA: DISTRIBUZIONI DI AUTOVALORI NON UNIFORMI

Pb: $\{ \lambda_i(A) \} = \sum_{i=1}^n \lambda_i \in [a, b]$ intervalli $U \{ \lambda_1' \}$ $b < \lambda_1'$



Il punto è che una volta sistemato il problema su $[a, b]$ con un'iterazione sistemiamo anche λ_1' .

Abbiamo dal Teorema

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq \min_{p \in P_k} \max_{\lambda \in \Sigma_A} |p(\lambda)|.$$

Prendiamo $k = k^*(a, b, \epsilon) + 1$ e un polinomio p particolare (che non è p^*)

$$| \leq \max_{\lambda \in \Sigma_A} |p(\lambda)|. \quad (*)$$

Come scelgo p ? Scelgo $p \in P_k = P_{k^*(a, b, \epsilon) + 1}$

$$p(\lambda) = \overset{\substack{\uparrow \\ \text{polinomio di} \\ \text{miglior approssima-} \\ \text{zione in } [a, b]}}{p^*(\lambda)} \cdot e(\lambda) \leftarrow e(\lambda) \text{ ad hoc: } e(\lambda) = 1 - \frac{\lambda}{\lambda_1'}$$

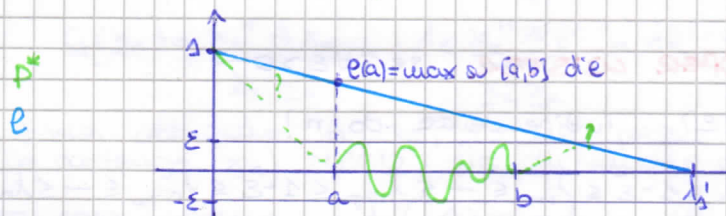
Osservo che $p(b) = p^*(b) \cdot e(b) = 1$ (ok!)

Abbiamo:

$$(*) \leq \max \left\{ \max_{\lambda \in [a, b]} |p(\lambda)|, |p(\lambda_1')| \right\} = \max_{\lambda \in [a, b]} |p(\lambda)|$$

\hookrightarrow nel continuo \parallel per costruzione

Facciamo un disegno per coprire la situazione:



$$0 < e(\lambda) \big|_{[a, b]} < e(a) = 1 - \frac{a}{\lambda_1'} < 1$$

Quindi abbiamo che:

$$\max_{\lambda \in [a, b]} |p(\lambda)| \leq \max_{\lambda \in [a, b]} |p^*(\lambda)| < \epsilon$$

\uparrow $\lambda \in [a, b]$
perché $e(\lambda) < 1$ in $[a, b]$

In conclusione ho fatto vedere che per aggiungere un autovalore fuori da $[a, b]$ non dobbiamo buttare via tutto, ma possiamo tenere conto di quello che abbiamo su $[a, b]$.

* 1. $\Sigma_A \subset [a,b] \cup (\bigcup_{i=1}^q \lambda_i')$ $b < \lambda_1'$ (outlier dx)

$\Rightarrow N(\varepsilon) = k^*(a,b,\varepsilon) + q$
iter

Quindi q outlier a dx
 $\Rightarrow q$ iterazioni in più

* 2. $\Sigma_A \subset [a,b] \cup (\bigcup_{i=1}^p \lambda_i)$ $\lambda_i < a$ (outlier a sx)

$\Rightarrow N(\varepsilon) = k^*(a,b,\varepsilon^*) + p$ con $\varepsilon^* = \varepsilon \prod_{i=1}^p \frac{\lambda_i}{b}$

Quindi nel caso degli outlier a sx cambia, infatti, pensiamo di avere uno solo, usiamo

$\rho(\lambda) = 1 - \frac{\lambda}{\lambda_1}$ su $[a,b]$ ma $|1 - \frac{\lambda}{\lambda_1}| < \frac{b}{\lambda_1}$ Non di 1! non possiamo metterlo via

Se facessimo i conti come prima troveremmo:

$\max_{\lambda \in \Sigma_A} |\rho(\lambda)| < \max_{\lambda \in [a,b]} |\rho^*(\lambda)| \cdot \prod_{i=1}^p \frac{b}{\lambda_i} < \varepsilon$ se $\max_{\lambda \in [a,b]} |\rho^*(\lambda)| < \varepsilon \prod_{i=1}^p \frac{\lambda_i}{b} = \varepsilon^*$

Dato che $\varepsilon^* < \varepsilon$, non ho solo p iterazioni in più perché arro' due
 $k^*(a,b,\varepsilon^*) > k^*(a,b,\varepsilon)$.

* 3. $\Sigma_A \subset [a,b] \cup (\bigcup_{i=1}^q \lambda_i) \cup (\bigcup_{i=1}^q \lambda_i')$

$\lambda_i < a$ outlier
 $\lambda_i' > b$

$\Rightarrow N(\varepsilon) = k^*(a,b,\tilde{\varepsilon}^*) + 2q$

con $\tilde{\varepsilon}^* = \varepsilon \prod_{i=1}^q \frac{\lambda_i}{\lambda_i'} \left(1 - \frac{\lambda_i}{\lambda_i'}\right)^{-2} < \varepsilon$
(BESTIO, ci interessa solo che $\tilde{\varepsilon} < \varepsilon$)

Siano $\{A_n\}$ e $\{P_n\}$ succ. di matrici hermitiane.

def: Dico che c'è **SPECTRAL EQUIVALENCE** tra $\{A_n\}$ e $\{P_n\}$ se

$\lambda_i(P_n^{-1}A_n) \in [d_1, d_2]$ $\forall i \forall n$

con d_1, d_2 costanti positive indipendenti dalla dimensione.

Si dice anche che $\{P_n\}$ è **PRECONDIZIONATORE OTTIMALE** per $\{A_n\}$.

def: Dico che $\{A_n\}$ ha un **WEAK CLUSTER** in p (valore) ^{es $p=1$} se $\forall \varepsilon > 0$, posto

$\gamma_{n,\lambda}(\varepsilon) = \# \{i: \lambda_i(A_n) \notin (p-\varepsilon, p+\varepsilon)\}$

ho $\lim_{n \rightarrow \infty} \frac{\gamma_{n,\lambda}(\varepsilon)}{n} = 0$. Ovvero ho pochi outlier rispetto alla dimensione

OSS: Se ho un cluster debole gli outlier possono aumentare ogni iterazione, ma in ogni caso la loro percentuale diminuisce (tende a zero, e' (n)).

def: Dico che ha un **PROPER CLUSTER** ^{CLUSTER FORTE} se $\forall \varepsilon > 0$

$\gamma_{n,\lambda}(\varepsilon) \leq C(\varepsilon)$ indipendente da n !

* 4. Σ_A t.c. $0 < \delta \leq \lambda_1 \leq \dots \leq \lambda_i < 1-\tilde{\varepsilon} \leq \lambda_{i+1} \leq \dots \leq \lambda_{n-j} \leq 1+\tilde{\varepsilon} \leq \lambda_{n-j+1} \leq \dots \leq \lambda_n$
(quindi ho un cluster a 1, i outlier sx e j outlier dx).

$\Rightarrow \frac{\|e_k\|_A}{\|e_{0,k}\|_A} \leq 2 \left(\frac{1+\tilde{\varepsilon}}{\delta} \right)^i \frac{1}{\varepsilon} \sim k-i-j$ $k > i+j = \# \text{ outlier}$

AUTOVALORI MATRICI HERMITIANE

28/11/2012
(15)

def: **QUOZIENTE DI RAYLEIGH**

$r_A(x) = \frac{x^H A x}{x^H x}$

$x \in \mathbb{C}^n \setminus \{0\}$
 $A \in \mathbb{C}^{n \times n}$ hermitiana

MEMO: 1. $S \subseteq \mathbb{C}^n$ sottosp $\leadsto \dim S^\perp = n - \dim S$

2. $S, T \subseteq \mathbb{C}^n$ sottosp $\leadsto \dim(S \cap T) \geq \max\{0, \dim S + \dim T - n\}$

3. $A \in \mathbb{C}^{m \times n}$ msn $S \subseteq \mathbb{C}^n$ sottosp

$$T = \{x \in \mathbb{C}^n : x = Ay \text{ } y \in S\} = \text{Im}_S(A)$$

$$\Rightarrow \dim T + \dim N = \dim S$$

$$N = \{x \in S : Ax = 0\} = \ker_S(A)$$

TEO: (COURANT-FISHER o minmax)

$A \in \mathbb{C}^{n \times n}$ Hermitiana con autovalori $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ allora

$$\lambda_{n-k+1}(A) = \min_{V_k} \max_{\substack{x \in V_k \\ x \neq 0}} r_A(x) \quad \forall k=1, \dots, n \quad V_k \subseteq \mathbb{C}^n \\ \dim V_k = k$$

$$\lambda_k(A) = \max_{V_k} \min_{\substack{x \in V_k \\ x \neq 0}} r_A(x)$$

dim: x_1, \dots, x_n autovettori ortonom. relativi agli autovalori λ_i

$$S = \langle x_k, \dots, x_n \rangle \rightarrow \dim S = n - k + 1$$

Ma dal pro de del MEMO ho

$$\dim(S \cap V_k) \geq \max\{0, (n-k+1) - k - n\} = 1$$

$$\Rightarrow S \cap V_k \neq \{0\}$$

Quindi possiamo prendere $x = \sum_{i=k}^n \alpha_i x_i \in S \cap V_k$. Abbiamo:

$$r_A(x) = \frac{x^H A x}{x^H x} = \frac{\left(\sum_{i=k}^n \bar{\alpha}_i x_i^H\right) A \left(\sum_{j=k}^n \alpha_j x_j\right)}{\sum_{i=k}^n \bar{\alpha}_i \alpha_i}$$

porto dentro la matrice e uso che $Ax_i = \lambda_i x_i$

poi uso anche $x_i \perp x_j$ se $i \neq j$ quindi rimane una somma sola

$$= \frac{\sum_{j=k}^n |\alpha_j|^2 \lambda_j x_j^H x_j}{\sum_{j=k}^n |\alpha_j|^2 x_j^H x_j} = \frac{\sum_{j=k}^n |\alpha_j|^2 \lambda_j}{\sum_{j=k}^n |\alpha_j|^2}$$

$$\leq \left(\max_{i=k-n} \lambda_i\right) \cdot \frac{\sum_{j=k}^n |\alpha_j|^2}{\sum_{j=k}^n |\alpha_j|^2} = \lambda_k$$

perché li avevo ordinati in ordine decrescente.

Quindi ho dimostrato che $r_A(x) \leq \lambda_k$, ma allora ovviamente

$$\min_{\substack{x \in V_k \\ x \neq 0}} r_A(x) \leq \lambda_k \quad (*)$$

Ora considero x_k e scelgo $V_k = \langle x_1, \dots, x_k \rangle$

$$r_A(x_k) = \frac{x_k^H A x_k}{x_k^H x_k} = \lambda_k \frac{x_k^H x_k}{x_k^H x_k} = \lambda_k$$

e quindi abbiamo finito, (perché posso scegliere V_k in modo che contenga x_k) DUBBIO: ma chi mi dice che su V_k il minimo è x_k

Ora consideriamo $\lambda_k(-A) = -\lambda_{n-k+1}(A)$

$$\begin{aligned} \max_{V_k} \left(\min_{\substack{x \in V_k \\ x \neq 0}} r_A(x) \right) &= \max_{V_k} \left(- \max_{\substack{x \in V_k \\ x \neq 0}} r_A(x) \right) = \\ &= \min_{V_k} (-f) = -\max_{V_k} f \\ &= -\min_{V_k} \max_{\substack{x \in V_k \\ x \neq 0}} (r_A(x)) \end{aligned}$$

Conseguenza:

$$\lambda_1 = \lambda_{\max}(A) = \max_{x \neq 0} r_A(x)$$

$V_n = \mathbb{C}^n$ è l'unico sottospazio di dimensione n ($k=n$) e la 1^a)

$$\lambda_n = \lambda_{\min}(A) = \min_{x \neq 0} r_A(x)$$

TEO: $A \in \mathbb{C}^{n \times n}$ Hermitiana, $U \in \mathbb{C}^{n \times n-1}$ t.c. $U^H U = I_{n-1}$, $B = U^H A U \in \mathbb{C}^{(n-1) \times (n-1)}$. Allora

compressione
di ordine
1

$$\lambda_1(A) \geq \lambda_1(B) \geq \lambda_2(A) \geq \dots \geq \lambda_{n-1}(A) \geq \lambda_{n-1}(B) \geq \lambda_n(A)$$

ovvero gli autovalori di B (che sono uno di meno) separano gli autovalori di A .

dimi: Per prima cosa mostro che $\lambda_{k-1}(B) \geq \lambda_k(A)$ per $k=2, \dots, n$

Uso la 2^a delle tesi del teorema di CF. So che $\exists z_k \in \mathbb{C}^n$, $\dim z_k = k$ tale che

$$\lambda_k = \min_{\substack{x \in z_k \\ x \neq 0}} r_A(x) \quad (\text{i.e. } z_k \text{ è lo spazio che realizza il max})$$

Scrivo $U = [u_1 | \dots | u_{n-1}]$ e pongo $S = \langle u_1, \dots, u_{n-1} \rangle \rightarrow \dim S = n-1$

$$\rightarrow \dim(z_k \cap S) \geq \max\{0, \dim z_k + \dim S - n\} = k-1 \geq 1$$

Quindi $\exists x \in z_k \cap S, x \neq 0$ t.c.

$$\lambda_k \leq \min_{\substack{x \in z_k \cap S \\ x \neq 0}} r_A(x)$$

Ora $\forall x \in S, \exists! y \in \mathbb{C}^{n-1}$ t.c. $x = Uy$ ($y = U^H x$)

Consideriamo $W = \{y \in \mathbb{C}^{n-1} : y = U^H x, x \in z_k \cap S\}$

$$\dim W + \underbrace{\dim \{x \in S \cap z_k : U^H x = 0\}}_{=0 \text{ per come è fatto } U} = \dim(S \cap z_k) \geq k-1$$

$$\rightarrow \dim W \geq k-1$$

Quindi abbiamo:

$$\lambda_k(A) \leq \min_{\substack{x \in z_k \cap S \\ x \neq 0}} r_A(x) = \min_{\substack{x \in z_k \cap S \\ x \neq 0}} \frac{x^H A x}{x^H x} = \min_{\substack{y \in W \\ y \neq 0}} \frac{y^H U^H A U y}{y^H U^H U y} = \min_{\substack{y \in W \\ y \neq 0}} \frac{y^H B y}{y^H y}$$

$$\leq \max_{V_{k-1}} \min_{\substack{y \in V_{k-1} \\ y \neq 0}} r_B(x) = \lambda_{k-1}(B)$$

Se W ha dimensione maggiore di uno, sta considerando tutti i suoi sottospazi, quindi il max è il più grande.

Per concludere applico quanto dimostrato a $-A$ e $-B$:

$$-\lambda_n(A) \geq -\lambda_{n-1}(A) \geq \dots \geq -\lambda_2(A)$$

$$-\lambda_{n-1}(B) \geq -\lambda_{n-2}(B) \geq \dots \geq -\lambda_1(B)$$

Ho che $\lambda_{k-1}(B) \geq \lambda_k(-A)$ $k=2, \dots, n$

$$\lambda_{(n-1)-(k-1)+1}(B) = \lambda_{n-k+1}(B)$$

$$\lambda_{n-k+1}(B)$$

Quindi eligendo $i = n-k+1$ ho trovato $\lambda_i(A) \geq \lambda_i(B)$ $i=1, \dots, n-1$

TEO: $A, B \in \mathbb{C}^{n \times n}$ Hermitiane con $B = A + \sigma \underline{u} \underline{u}^H$ correzione di rango 1
 $\underline{u} \in \mathbb{C}^n, \sigma \geq 0$

$$\Rightarrow \lambda_1(A) + \sigma \underline{u}^H \underline{u} \geq \lambda_1(B) \geq \lambda_2(A) \geq \dots \geq \lambda_{n-1}(A) \geq \lambda_n(B) \geq \lambda_n(A)$$

dim: se $\underline{u} = 0$ è ovvio dato che $A=B$. Sia $\underline{u} \neq 0$ vogliamo mostrare che

$$\lambda_i(B) \geq \lambda_i(A) \quad i=1 \dots n$$

Per il TEO CF (1) ho due per $k=1 \dots n$ $\exists Z_k \subseteq \mathbb{C}^n$ s.t. $\dim Z_k = k$ t.c.

$$\lambda_{n-k+1}(B) = \max_{\substack{x \in Z_k \\ x \neq 0}} r_B(x) = \max_{\substack{x \in Z_k \\ x \neq 0}} \left(\frac{x^H A x}{x^H x} + \sigma \frac{x^H \underline{u} \underline{u}^H x}{x^H x} \right)$$

$$\geq \max_{\substack{x \neq 0 \\ x \in Z_k}} \frac{x^H A x}{x^H x} \geq \min_k \max_{x \neq 0, x \in V_k} r_A(x) = \lambda_{n-k+1}(A) \quad \underline{dk}$$

Sempre dalla (1) del TEO CF ho che $\exists W_k \subseteq \mathbb{C}^n, \dim W_k = k$ t.c.

$$\lambda_{n-k+1}(A) = \max_{\substack{x \in W_k \\ x \neq 0}} r_A(x) \quad \text{i.e. } W_k \text{ è quello per cui vale } \underline{e}' = \text{qui } \otimes$$

Sia ora $S = \langle \underline{u} \rangle$ e consideriamo $T = W_k \cap S^\perp$ e prendiamo $x \in T$, ovviamente ho che $\underline{u}^H x = 0$. (dato che $x \in S^\perp$)

$$x^H B x = x^H A x + \sigma \frac{x^H \underline{u} \underline{u}^H x}{\underset{=0}{x^H x}} = x^H A x \quad \forall x \in T.$$

Allora abbiamo:

$$\lambda_{n-k+1}(A) \geq \max_{\substack{x \in T \\ x \neq 0}} r_A(x) = \max_{\substack{x \in T \\ x \neq 0}} r_B(x) \quad (**)$$

perché \uparrow
mi sto
restringendo
a T

Abbiamo due: $\dim S = 1 \rightarrow \dim S^\perp = n-1$ e quindi

$$\dim(T) = \dim(W_k \cap S^\perp) \geq \max\{0, \dim W_k + \dim S^\perp - n\} \\ = \max\{0, k + n - 1 - n\} = k - 1$$

Quindi da (**) posso dire che

$$\lambda_{n-k+1}(A) \geq \min_{V_{k-1}} \max_{\substack{x \neq 0 \\ x \in V_{k-1}}} r_B(x) = \lambda_{n-k+2}(B)$$

perché T ha $\dim \geq 0$ e quindi posso considerare tutti i suoi sspazi e poi \geq perché ci metto il min

Ora mi rimangono da considerare i λ_1 .

$$\lambda_1(B) = \max_{\substack{x \neq 0 \\ \text{CF}}} r_B(x) = \max_{x \neq 0} \left[\frac{x^H A x}{x^H x} + \sigma \frac{|x^H \underline{u}|^2}{x^H x} \right] \\ \stackrel{\text{Cauchy-Schwarz}}{\leq} \max_{x \neq 0} r_A(x) + \sigma \underline{u}^H \underline{u} = \lambda_1(A) + \sigma \underline{u}^H \underline{u}$$

GENERALIZZAZIONE: correzione di rango 3 $B = A + \sum_{i=1}^3 \sigma_i \underline{u}_i \underline{u}_i^H$

o analogamente al teo precedente $B = U^H A U$ con $U \in \mathbb{C}^{n \times m}$ con $m \leq n-1$ (e le idee volgono uguali...)